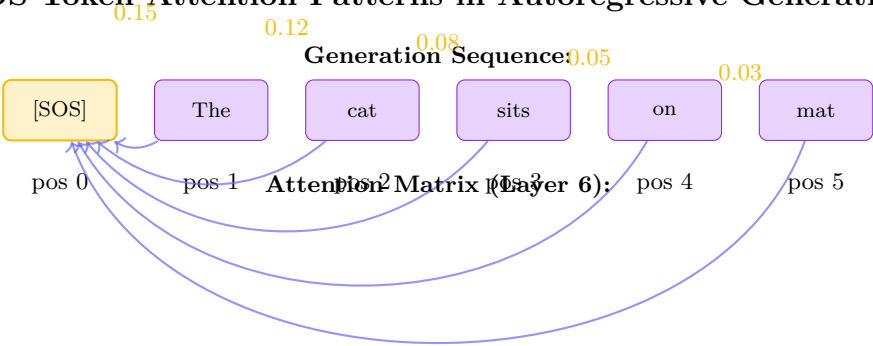


SOS Token Attention Patterns in Autoregressive Generation



Key Insights

[SOS] Token

Attends to SOS for global context

Persistent conditioning signal

Causal masking prevents future information leakage

Bidirectional attention flow

Generation Process: Each token uses SOS as initialization context

